



# INSTITUTE OF AERONAUTICAL ENGINEERING

(Autonomous)

Dundigal, Hyderabad - 500 043

## INFORMATION TECHNOLOGY

### TUTORIAL QUESTION BANK

<b>Course Title</b>	<b>DATA WAREHOUSING AND DATA MINING</b>
<b>Course Code</b>	<b>A60520</b>
<b>Regulation</b>	<b>R15 – JNTUH</b>
<b>Class</b>	<b>III B. Tech II Semester</b>
<b>Year</b>	<b>2017 – 2018</b>
<b>Course Faculty</b>	<b>Mr. Ch Suresh Kumar Raju, Assistant Professor, IT</b>

### OBJECTIVES

To meet the challenge of ensuring excellence in engineering education, the issue of quality needs to be addressed, debated and taken forward in a systematic manner. Accreditation is the principal means of quality assurance in higher education. The major emphasis of accreditation process is to measure the outcomes of the program that is being accredited.

In line with this, Faculty of Institute of Aeronautical Engineering, Hyderabad has taken a lead in incorporating philosophy of outcome based education in the process of problem solving and career development. So, all students of the institute should understand the depth and approach of course to be taught through this question bank, which will enhance learner's learning process.

<b>S. No</b>	<b>Question</b>	<b>Blooms Taxonomy Level</b>	<b>Course Outcome</b>
<b>UNIT - I</b>			
<b>PART – A (Short Answer Questions)</b>			
1	Define online analytical processing?	Remember	3
2	List the key features of data warehouse?	Understand	3
3	Define data mart?	Remember	3
4	Define enterprise warehouse?	Remember	3
5	Define virtual warehouse?	Remember	4
6	List the metadata repository?	Understand	4
7	List the various multidimensional models?	Understand	4
8	Explain about the star schema?	Understand	4
9	Explain the snowflake schema?	Understand	4
10	Define about the fact constellation model?	Remember	5
11	Name the OLAP operations?	Understand	1
12	Express what is slice and dice operation?	Understand	1
13	Define Pivot operation?	Remember	1
14	Distinguish between the OLAP Systems and Statistical databases?	Understand	1
15	State the various views of data warehouse design?	Understand	1
16	Define Relational OLAP(ROLAP) server?	Remember	3
17	Explain Multidimensional OLAP(MOLAP) server?	Understand	3
18	State what is Hybrid OLAP(HOLAP) server?	Understand	4

S. No	Question	Blooms Taxonomy Level	Course Outcome
19	Define Data warehouse?	Remember	3
20	Define the use of concept hierarchy?	Remember	4
<b>Part - B (Long Answer Questions)</b>			
1	Differentiate operational database systems and data warehousing?	Understand	8
2	Discuss briefly about the multidimensional data models?	Understand	7
3	Explain with an example the different schemas for multidimensional databases?	Understand	1
4	Describe the three-tier data warehousing architecture?	Remember	10
5	Discuss the efficient processing of OLAP queries?	Understand	5
6	Explain the data warehouse applications?	Understand	4
7	Explain the architecture for on-line analytical mining?	Understand	3
8	Describe the common techniques are used in ROLAP and MOLAP?	Remember	10
9	Describe the complex aggregation at multiple granularity?	Remember	4
10	Explain about the concept description? And what are the differences between concept description in large databases and OLAP?	Understand	3
11	Discuss about Metadata Repository?	Understand	5
12	Compare the schemas for the multidimensional data models?	Understand	4
13	Explain about the data warehouse implementation with an example?	Understand	5
14	Discuss about types of OLAP Servers?	Understand	5
15	Explain OLAP operations in the Multidimensional Data Model?	Understand	3
16	Compare Enterprise warehouse, data mart, virtual warehouse?	Understand	4
17	Compare Data cleaning, data transformation?	Understand	4
18	Explain what are the differences between the three main types of data warehouse usage: information processing, analytical processing and data mining? Discuss the motivation behind OLAP mining(OLAM)?	Understand	7
19	Explain a data warehouse can be modeled by either a star schema or a snowflake schema. Briefly describe the similarities and the differences of the two models, and then analyze their advantages and disadvantages with regard to one another?	Understand	3
20	Explain Indexing OLAP Data?	Understand	3
<b>Part - C (Problem Solving and Critical Thinking Questions)</b>			
1	Analyze that a data warehouse consists of the three dimensions time, doctor and patient, and the two measures count and charge, where charge is the fee that a doctor charges a patient for a visit. (a) Enumerate three classes of schemas classes of schemas that are popularly used for modeling data warehouses. (b) Draw a schema diagram for the above data warehouse using one of the schema classes listed in (a). (c) Starting with the base cuboid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2004? (d) To obtain the same list, write an SQL query assuming the data is stored in a relational database with the schema fee (day, month, year, doctor, hospital, patient, count, charge).	Understand	3
2	State why, for the integration of multiple heterogeneous information sources, many companies in industry prefer the update-driven approach (which constructs and uses data warehouses), rather than the query-driven approach (which applies wrappers and integrators). Describe situations where the query-driven approach is preferable over the update-driven approach.	Remember	3
3	Suppose that a data warehouse for Big University consists of the following four dimensions: student, course, semester, and instructor, and two measures count and average grade. When at the lowest conceptual level	Understand	3

	<p>(e.g., for a given student, course, semester, and instructor combination), the average grade measure stores the actual course grade of the student. At higher combination.</p> <p>Compute the number of cuboids</p> <p>(a) Draw a snowflake schema diagram for the data warehouse.</p> <p>(b) Starting with the base cuboid [student, course, semester, instructor], what specific OLAP operations (e.g., roll-up from semester to year ) Should one perform in order to list the average grade of CS courses for each Big University student.</p> <p>(c) If each dimension has five levels (including all), such as “student &lt; major &lt; status &lt; university &lt; all”, how many cuboids will this cube contain(including the base and apex cuboids)?</p>		
4	<p>Suppose that a data warehouse consists of the four dimensions, date, spectator, location, and game, and the two measures, count and charge, where charge is the fare that a spectator pays when watching a game on a given date. Spectators may be students, adults, or seniors, with each category having its own charge rate.</p> <p>Write the following</p> <p>(a) Draw a star schema diagram for the data warehouse.</p> <p>(b) Starting with the base cuboid [ date, spectator, location, game], what specific OLAP operations should one perform in order to list the total charge paid by student spectators at GM Place in 2004?</p> <p>(c) Bitmap indexing is useful in data warehousing. Taking this cube as an example, briefly discuss advantages and problems of using a bitmap index structure.</p>	Understand	3
5	<p>Design a data warehouse for a regional weather bureau. The weather bureau has about 1,000 probes, which are scattered throughout various land and ocean locations in the region to collect basic weather data, including air pressure, temperature, and precipitation at each hour. All data are sent to the central station, which has collected such data for over 10 years. Your design should facilitate efficient querying and on-line analytical processing, and derive general weather patterns in multidimensional space.</p>	Understand	8
6	<p>Explain the computation of measures in a data cube:</p> <p>(a) Enumerate three categories of measures, based on the kind of aggregate functions used in computing a data cube.</p> <p>(b) For a data cube with the three dimensions time, location, and item, which category does the function variance belong to? Describe how to compute it if the cube is partitioned into many chunks. Hint: The formula for computing variance is</p> $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}_i)^2$ <p>Where <math>\bar{x}_i</math> is the average of N <math>x_i</math>'s.</p> <p>(c) Suppose the function is “top 10 sales”. Discuss how to efficiently compute this measure in a data cube.</p>	Understand	8
7	<p>Suppose that we need to record three measures in a data cube: min, average, and median. Design an efficient computation and storage method for each measure given that the cube allows data to be deleted incrementally (i.e., in small portions at a time) from the cube.</p>	Understand	8
8	<p>Observe that a data warehouse contains 20 dimensions, each with about five levels of granularity.</p> <p>(a) Users are mainly interested in four particular dimensions, each having three frequently accessed levels for rolling up and drilling down. How would you design a data cube structure to efficiently support this preference?</p> <p>(b) At times, a user may want to drill through the cube, down to the raw data</p>	Remember	8

	for one or two particular dimensions. How would you support this feature?		
9	Observe A data cube, C, has n dimensions, and each dimension has exactly p distinct values in the base cuboid. Assume that there are no concept hierarchies associated with the dimensions. (a) What is the maximum number of cells possible in the base cuboid? (b) What is the minimum number of cells possible in the base cuboid? (c) What is the maximum number of cells possible (including both base cells and aggregate cells) in the data cube, C? (d) What is the minimum number of cells possible in the data cube, C?	Remember	8
10	Observe A popular data warehouse implementation is to construct a multidimensional database, known as a data cube. Unfortunately, this may often generate a huge, yet very sparse multidimensional matrix. Present an example illustrating such a huge and sparse data cube.	Remember	8

## UNIT - II

### Part – A (Short Answer Questions)

1	Define data mining?	Remember	1
2	Explain the definition of data warehouse?	Understand	1
3	Distinguish between data mining and data warehouse?	Understand	2
4	Identify any three functionality of data mining?	Remember	3
5	Interpret major issues in data mining?	Understand	1
6	Name the steps in the process of Remember discovery?	Remember	1
7	Discuss relational databases?	Understand	1
8	State object –oriented Databases?	Understand	1
9	Explain the spatial databases?	Understand	2
10	Contrast heterogeneous databases and legacy databases?	Understand	2
11	Differentiate classification and Prediction?	Understand	2
12	Describe transactional data bases?	Remember	2
13	List the types of data that can be mined?	Remember	3
14	Define data cube?	Remember	3
15	Define multidimensional data mining?	Remember	3
16	Define data characterization?	Remember	3
17	Express what is a decision tree?	Understand	3
18	Explain the outlier analysis?	Understand	3
19	Name the steps involved in data preprocessing?	Understand	3
20	Interpret the dimensionality reduction?	Understand	3

### Part - B (Long Answer Questions)

1	Describe data mining? In your answer, address the following: a)Is it another hype? b)Is it a simple transformation of Technology developed from databases, statistics, and machine learning? c)Explain how the evolutions of database technology lead to data mining? d)Describe the steps involved in data mining when viewed as a process of Remember discovery.	Understand	5
2	Distinguish between the data warehouse and databases? How they are similar?	Remember	4
3	Explain the difference between discrimination and classification? Between characterization and clustering? Between classification and prediction? For each of these pairs of tasks, how are they similar?	Understand	4
4	Describe three challenges to data mining regarding data mining methodology and user interaction issues?	Remember	4
5	Distinguish between the data warehouses and data mining?	Remember	5
6	Discuss briefly about the data smoothing techniques?	Understand	4
7	Explain Data Integration and Transformation?	Understand	6
8	Describe the various data reduction techniques?	Understand	8

9	Define data cleaning? Express the different techniques for handling missing values?	Remember	5
10	Differentiate between descriptive and predictive data mining?	Understand	3
11	Explain data mining as a step in the process of Remember discovery?	Understand	6
12	Describe briefly Discretization and concept hierarchy generation for numerical data?	Remember	6
13	Discuss about the concept hierarchy generation for categorical data?	Understand	7
14	List and describe the five primitives for specifying a data mining task?	Understand	4
15	Discuss issues to consider during data integration?	Understand	1
16	Describe the following advanced database systems and applications: object- relational databases, spatial databases, text databases, multimedia databases, stream data, the World Wide Web.	Remember	5
17	Describe why concept hierarchies are useful in data mining.	Remember	5
18	Describe the differences between the following approaches for the integration of a data mining system with a database or data warehouse system: no coupling, loose coupling, semi tight coupling, and tight coupling. State which approach you think is the most popular, and why	Remember	1
19	Explain Data quality can be assessed in terms of accuracy, completeness, and consistency. Propose two other dimensions of data quality.	Understand	1
20	Understand the two methods below to normalize the following group of data: 200, 300, 400, 600, 1000  (a) min-max normalization by setting min = 0 and max = 1 (b) z-score normalization	Understand	8
<b>Part – C (Problem Solving and Critical Thinking)</b>			
1	Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52,70. Compute the following: (a) Mean of the data? Median? (b) mode of the data? Comment on the data's modality (i.e. bimodal, trimodal etc.). (c) midrange of the data?	Understand	1
2	Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19,20,20,21,22,22,25,2525, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52,70. Compute the following: (a) Can you find (roughly) the first quartile (Q1) and the third quartile (Q3) of the data? (b) Give the five-number summary of the data. (c) Show a box plot of the data. (d) How is a quantile-quantile plot different from a quantile plot?	Understand	1
3	Use the data for age given above answer the following. (a) Use smoothing by bin means to smooth the above data, using a bin depth of 3. Illustrate your steps. Comment on the effect of this technique for the given data (b) How might you determine outliers in the data? (c) What other methods are there for data smoothing?	Understand	1
4	Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result age 23 23 27 27 39 41 47 49 50 %fat 9.5 26.5 7.8 17.8 31.4 25.9 27.4 27.2 31.2 age 52 54 54 56 57 58 58 60 61	Remember	9

	%fat 34.6 42.5 28.8 33.4 30.2 34.1 32.9 41.2 35.7 Examine the following (a) the mean, median and standard deviation of age and %fat. (b) Draw the box plots for age and %fat. (c) Draw a scatter plot and a q-q plot based on these two variables.		
5	Write an example where data mining is crucial to the success of a business. What data mining functions does this business need? Can they be performed alternatively by data query processing or simple statistical analysis?	Understand	6
6	Suppose your task as a software engineer at Big University is to design a data mining system to examine the university course database, which contains the following information: the name, address, and status (e.g., undergraduate or graduate) of each student, the courses taken, and the cumulative grade point average (GPA). Describe the architecture you would choose. What is the purpose of each component of this architecture?	Understand	8
7	Outliers are often discarded as noise. However, one person's garbage could be another's treasure. For example, exceptions in credit card transactions can help us detect the fraudulent use of credit cards. Taking fraudulence detection as an example, Write two methods that can be used to detect outliers and discuss which one is more reliable.	Understand	8
9	Examine the following consider the following data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70. (a) Use min-max normalization to transform the value 35 for age on to the range [0.0, 1.0]. (b) Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years. (c) Use normalization by decimal scaling to transform the value 35 for age. (d) Comment on which method you would prefer to use for the given data, giving reasons as to why.	Remember	8
10	Suppose a group of 12 sales price records has been sorted as follows : 5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215 Examine the following methods by partition them into three bins (a) equal-frequency (equi-depth) partitioning (b) equal-width partitioning (c) clustering	Remember	8

### UNIT-III

#### Part - A (Short Answer Questions) (MID-I)

1.	Define association rule?	Remember	3
2.	Define item set?	Remember	3
3.	Define frequent item sets?	Understand	3
4.	List the measures of association rules?	Understand	3
5.	Give the types of association rules?	Understand	1
6.	Give the principle of APRIORI algorithm?	Understand	6
7.	Specify the problem definition for association rules?	Remember	7
8.	What is support and minimum support?	Understand	4
9.	What is confidence and minimum confidence for strong association rule?	Understand	1
10.	Name the steps in association rule mining?	Remember	2
11.	Explain the two kinds of closure checking?	Understand	3
12.	Describe the five categories of pattern mining constraints?	Understand	1
13.	List the techniques of efficiency of Apriori algorithm?	Remember	1
14.	List the drawbacks of Apriori technique?	Understand	2
15.	Give examples for frequent item sets?	Understand	3

16	Name the pruning strategies in mining closed frequent item sets?	Understand	4
17	Explain the join step?	Understand	5
18	Describe the prune step?	Remember	1
19	State how can we mine closed frequent item sets?	Understand	6
20	Name the pruning strategies of closed frequent item sets?	Remember	7
21	Interpret the rule of support for item sets?	Understand	2
22	Explain the two kinds of closure checking?	Understand	3
23	List the techniques to improve the efficiency of Apriori algorithm?	Understand	2
24	Write the procedure to find association rule from given frequent item sets?	Understand	3

**(MID-II)**

1	Define Prediction?	Understand	7
2	Describe how is Euclidean distance measured?	Understand	6
3	Define Inductive logic programming?	Understand	5
4	Describe KDD?	Understand	5
5	Define FP-Growth?	Understand	5
6	Explain Minimum steps to construct FP tree for minimum 5 items?	Remember	7
7	Define item sets?	Understand	6
8	Define closed item sets?	Understand	6
9	List two techniques of efficiency of Apriori algorithm?	Understand	6
10	List the pruning strategies in mining closed frequent item sets?	Understand	6
11	Define substructure of a structural pattern?	Understand	7
12	Define frequent patterns?	Understand	5
13	Discuss the FP-growth algorithm?	Understand	5
14	Define Support?	Understand	5
15	Define Confidence?	Understand	6
16	Explain the join step?	Understand	7
17	Describe the prune step?	Remember	6
18	State how can we mine closed frequent item sets?	Remember	6
19	Discuss Sub-item set pruning?	Remember	6
20	Explain the two kinds of closure checking?	Remember	5
21	Describe two categories of pattern mining constraints?	Understand	5
22	Discuss the purpose of node-links in FP – tree?	Understand	5
23	Explain tree pruning?	Understand	5
24	Define the decision tree?		

**Part – B (Long Answer Questions)**

1	Define the terms frequent item sets, closed item sets and association rules?	Remember	9
2	Discuss which algorithm is an influential algorithm for mining frequent Item sets for boolean association rules? Explain with an example?	Understand	8
3	Describe the different techniques to improve the efficiency of Apriori? Explain?	Remember	4
4	Discuss the FP-growth algorithm? Explain with an example?	Understand	6
5	Explain how to mine the frequent item sets using vertical data format?	Understand	5
6	Discuss about mining multilevel association rules from transaction databases in detail?	Understand	4
7	Explain how to mine the multidimensional association rules from relational databases and data warehouses?		
8	Describe briefly about the different correlation measures in association analysis?		
9	Discuss about constraint-based association mining?	Understand	2
10	Explain the Apriori algorithm with example?	Understand	4
11	Discuss the generating association rules from frequent item sets.	Understand	4
	Discuss about mining multilevel association rules from transaction databases in detail?		
13	Describe multidimensional association rules using static Discretization?	Remember	4
14	Explain what are additional rule constraints to guide mining?	Understand	4

	Explain, how can we tell which strong association rules are really interesting? Explain with an example?														
16	Describe about the correlation analysis using Chi-square?	Remember	4												
17	<p>Understand the following rules on a database has five transactions. Let min sup = 60% and min con f = 80%.</p> <table border="1" style="margin-left: 40px;"> <thead> <tr> <th>TID</th> <th>items bought</th> </tr> </thead> <tbody> <tr> <td>T100</td> <td>{M, O, N, K, E, Y}</td> </tr> <tr> <td>T200</td> <td>{D, O, N, K, E, Y}</td> </tr> <tr> <td>T300</td> <td>{M, A, K, E}</td> </tr> <tr> <td>T400</td> <td>{M, U, C, K, Y}</td> </tr> <tr> <td>T500</td> <td>{C, O, O, K, I, E}</td> </tr> </tbody> </table> <p>(a) Find all frequent item sets using Apriori .  (b) List all of the strong association rules (with support s and confidence  (c) matching the following meta rule, where X is a variable representing customers, and item denotes variables representing items (e.g., “A”, “B”, etc.):  <math>\forall x \in \text{transaction}, \text{buys}(X, \text{item1}) \wedge \text{buys}(X, \text{item2}) \Rightarrow \text{buys}(X, \text{item3}) [s, c]</math></p>	TID	items bought	T100	{M, O, N, K, E, Y}	T200	{D, O, N, K, E, Y}	T300	{M, A, K, E}	T400	{M, U, C, K, Y}	T500	{C, O, O, K, I, E}	Understand	8
TID	items bought														
T100	{M, O, N, K, E, Y}														
T200	{D, O, N, K, E, Y}														
T300	{M, A, K, E}														
T400	{M, U, C, K, Y}														
T500	{C, O, O, K, I, E}														
18	Describe about the Mining closed Frequent Item set	Remember	8												
19	Write a short example to show that items in a strong association rule may actually be negatively correlated.	Understand	3												
20	Explain Association rule mining often generates a large number of rules. Discuss effective methods that can be used to reduce the number of rules generated while still preserving most of the interesting rules.	Understand	3												
<b>Part – C (Problem Solving and Critical Thinking Questions)</b>															
1	<p>The Apriori algorithm uses prior Remember of subset support properties. Analyze</p> <p>(a) That all nonempty subsets of a frequent item set must also be frequent.  (b) The support of any nonempty subset <math>s_0</math> of item set <math>s</math> must be at least as great as the support of <math>s</math>.  (c) Given frequent item set <math>l</math> and subset <math>s</math> of <math>l</math>, prove that the confidence of the rule “<math>s_0 \Rightarrow (l - s_0)</math>” cannot be more than the confidence of “<math>s \Rightarrow (l - s)</math>”, where <math>s_0</math> is a subset of <math>s</math>.  (d) A partitioning variation of Apriori subdivides the transactions of a database <math>D</math> into <math>n</math> non overlapping partitions. Prove that any item set that is frequent in <math>D</math> must be frequent in at least one partition of <math>D</math>.</p>	Understand	3												
2	<p>Implement three frequent item set mining algorithms introduced in this chapter :</p> <p>(1) Apriori [AS94], (2) FP-growth [HPY00], and (3) ECLAT [Zak00] (mining using vertical data format), using a programming language that you are familiar with, such as C++ or Java.  Compare the performance of each algorithm with various kinds of large data set.  Write a report to analyze the situations (such as data size, data distribution, minimal support threshold setting, and pattern density) where one algorithm may perform better than the others, and state why.</p>	Understand	6												
3	<p>Suppose that a large store has a transaction database that is distributed among four locations. Transactions in each component database have the same format , namely <math>T_j : \{i_1, \dots, i_m\}</math>, where <math>T_j</math> is a transaction identifier, and <math>i_k (1 \leq k \leq m)</math> is the identifier of an item purchased in the transaction. Construct an efficient algorithm to mine global association rules (without</p>	Understand	3												



	considering multilevel associations). You may present your algorithm in the form of an outline. Your algorithm should not require shipping all of the data to one site and should not cause excessive network communication overhead.																		
4	Suppose that frequent item sets are saved for a large transaction database, DB. Illustrate how to efficiently mine the (global) association rules under the same minimum support threshold if a set of new transactions, denoted as $\Delta DB$ , is (incrementally) added in?	Understand	3																
5	Most frequent pattern mining algorithms consider only distinct items in a transaction. However, multiple occurrences of an item in the same shopping basket, such as four cakes and three jugs of milk, can be important in transaction data analysis. Analyze how can one mine frequent item sets efficiently considering multiple occurrences of items? Propose modifications to the well-known algorithms, such as Apriori and FP-growth, to adapt to such a situation.	Understand	1																
6	<p>A database has five transactions. Let <math>\text{min sup} = 60\%</math> and <math>\text{min con f} = 80\%</math>.</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>TID</th> <th>items bought</th> </tr> </thead> <tbody> <tr> <td>T100</td> <td>{M, O, N, K, E, Y}</td> </tr> <tr> <td>T200</td> <td>{D, O, N, K, E, Y}</td> </tr> <tr> <td>T300</td> <td>{M, A, K, E}</td> </tr> <tr> <td>T400</td> <td>{M, U, C, K, Y}</td> </tr> <tr> <td>T500</td> <td>{C, O, O, K, I, E}</td> </tr> </tbody> </table> <p>Examine the following</p> <p>(d) Find all frequent item sets using FP-growth.</p> <p>(e) List all of the strong association rules (with support <math>s</math> and confidence <math>c</math>) matching the following meta rule, where <math>X</math> is a variable representing customers, and item <math>i</math> denotes variables representing items (e.g., "A", "B", etc.):</p> $\forall x \in \text{transaction}, \text{buys}(X, \text{item1}) \wedge \text{buys}(X, \text{item2}) \Rightarrow \text{buys}(X, \text{item3})$ <p>[<math>s, c</math>]</p>	TID	items bought	T100	{M, O, N, K, E, Y}	T200	{D, O, N, K, E, Y}	T300	{M, A, K, E}	T400	{M, U, C, K, Y}	T500	{C, O, O, K, I, E}	Remember	8				
TID	items bought																		
T100	{M, O, N, K, E, Y}																		
T200	{D, O, N, K, E, Y}																		
T300	{M, A, K, E}																		
T400	{M, U, C, K, Y}																		
T500	{C, O, O, K, I, E}																		
7	<p>The following contingency table summarizes supermarket transaction data, where hot dogs refers to the transactions containing hot dogs, hot dogs refers to the transactions that do not contain hot dogs, hamburgers refers to the transactions containing hamburgers, and hamburgers refers to the transactions that do not contain hamburgers</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>hot dogs</th> <th>hot dogs</th> <th>row</th> </tr> </thead> <tbody> <tr> <td>hamburgers</td> <td>2,000</td> <td>500</td> <td>2,500</td> </tr> <tr> <td>hamburgers</td> <td>1,000</td> <td>1,500</td> <td>2,500</td> </tr> <tr> <td>total</td> <td>3,000</td> <td>2,000</td> <td>5,000</td> </tr> </tbody> </table> <p>Observe that the association rule "hot dogs <math>\Rightarrow</math> hamburgers" is mined. Given a minimum support threshold of 25% and a minimum confidence threshold of 50%, is this association rule strong? Based on the given data, is the purchase of hot dogs independent of the purchase of hamburgers? If not, what kind of correlation relationship exists between the two</p>		hot dogs	hot dogs	row	hamburgers	2,000	500	2,500	hamburgers	1,000	1,500	2,500	total	3,000	2,000	5,000	Remember	8
	hot dogs	hot dogs	row																
hamburgers	2,000	500	2,500																
hamburgers	1,000	1,500	2,500																
total	3,000	2,000	5,000																
8	Sequential patterns can be mined in methods similar to the mining of	Understand	8																

	association rules. Design an efficient algorithm to mine multilevel sequential patterns from a transaction database. An example of such a pattern is the following: “A customer who buys a PC will buy Microsoft software within three months,” on which one may drill down to find a more refined version of the pattern, such as “A customer who buys a Pentium PC will buy Microsoft Office within three months.”		
9	The price of each item in a store is nonnegative. The store manager is only interested in rules of the form: “one free item may trigger \$200 total purchases in the same transaction.” Describe how to mine such rules efficiently.	Remember	8
10	The price of each item in a store is nonnegative. For each of the following cases, identify the kinds of constraint they represent and briefly discuss how to mine such association rules efficiently. (a) Containing at least one Nintendo game (b) Containing items the sum of whose prices is less than \$150 (c) Containing one free item and other items the sum of whose prices is at least \$200 (d) Where the average price of all the items is between \$100 and \$500	Understand	8

#### UNIT-IV

##### Part – A (Short Answer Questions)

1	State classification?	Understand	5
2	Define regression analysis?	Remember	4
3	Name the steps in data classification?	Understand	4
4	Define training tuple?	Remember	4
6	Describe accuracy of a classifier?	Remember	5
7	Differentiate supervised learning and unsupervised learning?	Understand	4
8	Define the decision tree?	Understand	5
9	Define information gain?	Remember	5
10	State gain ratio?	Understand	4
11	State Gini index?	Understand	4
12	Explain tree pruning?	Understand	4
14	Define the construction of naïve Bayesian classification?	Understand	5
15	Explain the IF-THEN rules for classification?	Understand	5
16	Explain Decision Tree Induction?	Understand	5
17	List the Attribute Selection Measures?	Remember	5
18	Define Bayes’ Theorem?	Understand	5
19	Define Naïve Bayesian Classification?	Remember	5
20	Explain K-Nearest-Neighbor Classifiers?	Understand	4

##### Part – B (Long Answer Questions)

1	Explain about the classification and prediction? Example with an example?		
2	Discuss about basic decision tree induction algorithm?	Understand	4
3	Explain briefly various measures associated with attribute selection?	Understand	3
4	Summarize how does tree pruning work? What are some enhancements to basic decision tree induction?	Understand	4
5	Explain how scalable is decision tree induction? Explain?	Understand	3
6	Describe the working procedures of simple Bayesian classifier?	Remember	4
7	Explain Bayesian Belief Networks?	Understand	5
8	Discuss about k-nearest neighbor classifier and case-based reasoning?	Understand	1
9	Explain about classifier accuracy? Explain the process of measuring the accuracy of a classifier?	Understand	2
10	Describe any ideas can be applied to any association rule mining be applied to classification?	Remember	3
11	Explain briefly about the Navie Bayesian Classification?	Remember	3
12	Explain about the major issues regarding classifications and predictions?	Understand	5
13	Differentiate classification and prediction methods?	Understand	5

14	Explain briefly various measures associated with attribute selection?	Understand	5
15	Explain training of Bayesian belief networks?	Understand	5
16	Explain how tree pruning useful in decision tree induction? What is a drawback of using a separate set of tuples to evaluate pruning?	Understand	7
17	Explain for a given a decision tree, you have the option of (a) converting the decision tree to rules and then pruning the resulting rules, or (b) pruning the decision tree and then converting the pruned tree to rules. What advantage does (a) have over (b)?	Understand	7
18	Compare the advantages and disadvantages of eager classification (e.g., decision tree, Bayesian, neural network) versus lazy classification (e.g., k-nearest neighbor, case-based reasoning).	Understand	7
19	Write an algorithm for k-nearest-neighbor classification given k and n, the number of attributes describing each tuple.	Understand	7
20	Describe each of the following clustering algorithms in terms of the following criteria: (i) shapes of clusters that can be determined; (ii) input parameters that must be specified; and (iii) limitations. (a)k-means (b)k-medoids	Remember	3
<b>Part – C (Problem Solving and Critical Thinking Questions)</b>			
1	Illustrate why is tree pruning useful in decision tree induction? Explain the drawback of using a separate set of tuples to evaluate pruning?	Understand	5
2	Given a decision tree, you have the option of (a) converting the decision tree to rules and then pruning the resulting rules, or (b) pruning the decision tree and then converting the pruned tree to rules. Explain advantage does (a) have over (b)?	Understand	4
3	Outline the major ideas of naïve Bayesian classification. Explain why is naïve Bayesian classification called “naïve”?	Understand	4
4	Design an efficient method that performs effective naïve Bayesian classification over an infinite data stream (i.e., you can scan the data stream only once). If we wanted to discover the evolution of such classification schemes (e.g., comparing the classification scheme at this moment with earlier schemes, such as one from a week ago),Construct modified design would you suggest?	Understand	4
5	The support vector machine (SVM) is a highly accurate classification method. However, SVM classifiers suffer from slow processing when training with a large set of data tuples. Explain how to overcome this difficulty and develop a scalable SVM algorithm for efficient SVM classification in large datasets.	Understand	5
6	It is important to calculate the worst-case computational complexity of the decision tree algorithm. Given data set D, the number of attributes n, and the number of training tuples  D , Show that the computational cost of growing a tree is at most $n \times  D  \times \log( D )$ .	Understand	7
7	Given a 5 GB data set with 50 attributes (each containing 100 distinct values) and 512 MB of main memory in your laptop, outline an efficient method that constructs decision trees in such large data sets. Justify your answer by rough calculation of your main memory usage.	Understand	7
8	What is associative classification? Why is associative classification able to achieve higher classification accuracy than a classical decision tree method? Explain how associative classification can be used for text document classification.	Understand	7
9	It is difficult to assess classification accuracy when individual data objects may belong to more than one class at a time. In such cases, Explain on what criteria you would use to compare different classifiers modeled after the same data.	Understand	7
10	Describe each of the following clustering algorithms in terms of the	Understand	3

	<p>following criteria: (i) shapes of clusters that can be determined; (ii) input parameters that must be specified; and (iii) limitations.</p> <p>(a) k-means  (b) k-medoids  (c) CLARA  (d) BIRCH  (e) ROCK  (f) Chameleon  (g) DBSCAN</p>		
<b>UNIT-V</b>			
<b>Part - A (Short Answer Questions)</b>			
1	Define Clustering?	Remember	5
2	Illustrate the meaning of cluster analysis?	Understand	5
3	Explain the fields in which clustering techniques are used?	Understand	5
4	List out the requirements of cluster analysis?	Remember	5
5	Express the different types of data used for cluster analysis?	Understand	5
6	State interval scaled variables?	Remember	5
7	Define Binary variables? And what are the two types of binary variables?	Remember	5
8	Define nominal, ordinal and ratio scaled variables?	Remember	6
9	Illustrate mean by partitioning method?	Understand	6
10	Define CLARA and CLARANS?	Remember	6
11	State hierarchical method?	Remember	6
12	Differentiate agglomerative and divisive hierarchical clustering?	Understand	6
13	State K-Means method?	Remember	6
14	Define Outlier Detection?	Remember	5
20	Define Chameleon method?	Remember	5
<b>Part - B (Long Answer Questions)</b>			
1	Discuss the various types of data in cluster analysis?	Understand	3
2	Explain the categories of major clustering methods?	Understand	3
3	Write algorithms for k-means and k-medoids? Explain?	Understand	3
4	Describe the different types of hierarchical methods?	Understand	3
5	Demonstrate about the following hierarchical methods a) BIRCH b) Chameleon	Understand	3
6	Explain about semi-supervised cluster analysis?	Understand	3
7	Explain about the outlier analysis?	Understand	3
8	Define the distance-based outlier? Illustrate the efficient algorithms for mining distance-based algorithm?	Remember	3
9	Explain about the Statistical-based outlier detection?	Understand	5
10	Describe about the distance-based outlier detection?	Remember	5
11	Discuss about the density-based outlier detection?	Understand	5
12	Demonstrate about the deviation-based outlier detection techniques?	Understand	5
13	Demonstrate about the BIRCH hierarchical methods?	Understand	5
14	Demonstrate about the ROCK(Robust Clustering using links) hierarchical methods?	Understand	5
15	Explain about the agglomerative and divisive hierarchical methods?	Understand	3
16	Demonstrate how to compute the dissimilarity between objects described by the following types of variables: (a) Numerical (interval-scaled) variables (b) Asymmetric binary variables (c) Categorical variables (d) Ratio-scaled variables (e) Non metric vector objects	Understand	6
17	Understand the following measurements for the variable age: 18, 22, 25, 42, 28, 43, 33, 35, 56, 28,	Understand	6

	standardize the variable by the following: (a) Compute the mean absolute deviation of age. (b) Compute the z-score for the first four measurements.		
18	Illustrate the strength and weakness of k-means in comparison with the k-medoids algorithm. Also, illustrate the strength and weakness of these schemes in comparison with a hierarchical clustering scheme (such as AGNES).	Understand	6
19	Explain why is outlier mining important? Briefly describe the different approaches behind statistical-based outlier detection, distanced-based outlier detection, density-based local outlier detection, and deviation-based outlier detection.	Understand	6
20	Understand the given following measurements for the variable age: 28, 32, 15, 42, 28, 43, 30, 32, 55, 26, standardize the variable by the following: (a) Compute the mean absolute deviation of age. (b) Compute the z-score for the first four measurements.	Understand	6
<b>Part – C (Problem Solving and Critical Thinking Questions)</b>			
1	Given the following measurements for the variable age: 48, 12, 25, 42, 28,43,33,35, 56, 28, standardize the variable by the following: Compute (a) The mean absolute deviation of age. (b)The z-score for the first four measurements.	Understand	5
2	Given two objects represented by the tuples (22, 1, 42, 10) and (20, 0, 36,8): Compute (a) The Euclidean distance between the two objects. (b)The Manhattan distance between the two objects. (c) The Minkowski distance between the two objects, using $p = 3$ .	Understand	5
3	Suppose that the data mining task is to cluster the following eight points(with (x, y) representing location) into three clusters. A1(2, 10), A2(2, 5), A3(8, 4), B1(5, 8),B2(7, 5), B3(6, 4), C1(1, 2), C2(4, 9). The distance function is Euclidean distance. Suppose initially we assign A1, B1, and C1 as the center of each cluster, respectively. Use the k-means algorithm to show only (a) The three cluster centers after the first round of execution and (b) The final three clusters	Understand	5
4	Explain why is it that BIRCH encounters difficulties in finding clusters of arbitrary shape but OPTICS does not? Can you propose some modifications to BIRCH to help it find clusters of arbitrary shape?	Understand	5
5	Clustering has been popularly recognized as an important data mining task with broad applications. Show one application example for each of the following cases: (a) An application that takes clustering as a major data mining function (b) An application that takes clustering as a preprocessing tool for data preparation for other data mining tasks	Understand	7
6	Clustering has been popularly recognized as an important data mining task with broad applications. Give example for each of the following cases: (a) An application that takes clustering as a major data mining function (b) An application that takes clustering as a preprocessing tool for data preparation for other data mining tasks	Understand	3
7	Data cubes and multidimensional databases contain categorical, ordinal, and numerical data in hierarchical or aggregate forms. Based on what you have learned about the clustering methods, Design a clustering method that finds clusters in large data cubes effectively and efficiently.	Understand	9
8	Human eyes are fast and effective at judging the quality of clustering methods for two- dimensional data. Design a data visualization method that may help humans visualize data clusters and judge the clustering quality for	Understand	3

	three-dimensional data? What about for even higher-dimensional data?		
9	Given the following measurements for the variable age: 29, 31, 25, 41, 27, 43, 33, 35, 56, 28, standardize the variable by the following: Compute (a) The mean absolute deviation of age. (b) The z-score for the first three measurements.	Understand	6
10	Given two objects represented by the tuples (21, 2, 41, 11) and (21, 1, 32, 6): Compute (a) The Euclidean distance between the two objects. (b) The Manhattan distance between the two objects. (c) The Minkowski distance between the two objects, using $p = 2$ .	Understand	8

Prepared by : Ch Suresh Kumar Raju, Assistant Professor, IT

**HOD,IT**